

Segmentation of the Poisson and negative binomial rate models: a penalized estimator

A. Cleynen¹ and E. Lebarbier¹

¹*AgroParisTech/ INRA MIA 518, 16 rue Claude Bernard, 75231 Paris Cedex 05, France.*

E-mail: alice.cleynen@agroparistech.fr

E-mail: emilie.lebarbier@agroparistech.fr

March 21, 2013

Abstract

We consider the segmentation problem of Poisson and negative binomial (i.e. overdispersed Poisson) rate distributions. In segmentation, an important issue remains the choice of the number of segments. To this end, we propose a penalized log-likelihood estimator where the penalty function is constructed in a non-asymptotic context following the works of L. Birgé and P. Massart. The resulting estimator is proved to satisfy an oracle inequality. The performances of our criterion is assessed using simulated and real datasets in the RNA-seq data analysis context.

Mathematics subject classification 2010: primary 62G05, 62G07; secondary 62P10

Keywords and phrases: Density estimation; Change-points detection; Count data (RNA-seq); Poisson and negative binomial distributions; Model selection.

Introduction

We consider a multiple change-point detection setting for count datasets, which can be written as follows: we observe a finite sequence $\{y_t\}_{t \in \{1, \dots, n\}}$ realisation of independent variables Y_t . These variables are supposed to be drawn from a probability distribution \mathcal{G} which depends on a set of parameters. Here two types of parameters are distinguished:

$$Y_t \sim \mathcal{G}(\theta_t, \phi) = s(t), \quad 1 \leq t \leq n,$$

where ϕ is a constant parameter while the θ s are point-specific. In many contexts, we might want to consider that the θ s are piece-wise constant and so subject to an unknown number $K - 1$ of abrupt changes (for instance with climatic or financial data). Thus, we want to assume the existence of partition of $\{1, \dots, n\}$ into K segments within which the observations follow the same distribution and between which observations have different distributions, i.e. θ is constant within a segment and differ from a segment to another.

A motivating example is sequencing data analysis. For instance, the output of RNA-seq experiments is the number of reads (i.e. short portions of the genome) which first position maps to each location of a genome of reference. Supposing that we dispose of such a sequence, we expect to observe a stationarity in the amount of reads falling in different areas of the genome: expressed genes, intronic regions, etc. We wish to localize those regions that are biologically significant. In our context, we consider for \mathcal{G} the Poisson and negative binomial distributions, adapted to RNA-seq experiment analysis [1].

Change-point detection problems are not new and many methods have been proposed in the literature. For count data-sets, [2] provide a detailed bibliography of methods in the particular case of the segmentation of the DNA sequences that includes Bayesian approaches, scan statistics, likelihood-ratio tests, binary segmentation and numerous other methods such as penalized contrast estimation procedures. In a Bayesian framework, [3] proposes to use an exact "ICL" criterion for the choice of K , while its approximation is computed in the constrained HMM approach of [4]. In this paper, we consider a penalized contrast estimation method which consists first, for every fixed K , in finding the best segmentation in K segments by minimizing the contrast over all the partitions with K segments, and then in selecting a convenient number of segments K by penalizing the contrast. Choosing the number of segments, i.e. choosing a "good" penalty, is a crucial issue and not so easy. The most basic examples of penalty are the Akaike Information Criterion (AIC [5]) and the Bayes Information Criterion (BIC [6]) but these criteria are not well adapted in the segmentation context and tend to overestimate the number of change-points (see [7, 8] for theoretical explanations). In this particular context, some modified versions of these criteria have been proposed. For instance, [8, 9] have proposed modified versions of the BIC criterion (shown to be consistent) in the segmentation of Gaussian processes and DNA sequences respectively. However, these criteria are based on asymptotic considerations. In the last years there has been an extensive literature influenced by [10, 11] introducing non-asymptotic model selection procedures, in the sense that the size of the models as well as the size of the list of models are allowed to be large when n is large. This penalized contrast procedure consists in selecting a model amongst a collection such that its performance is as close as possible to that of the best but unreachable model in terms of risk. This approach has been now considered in various function estimation contexts. In particular, [12] proposed a penalty for estimating the density of independent categorical variables in a least-squares framework, while [13, 14], or [15], focused on the estimation of the density of a Poisson process.

When the number of models is large, as in the case of an exhaustive search in segmentation problem, it can be shown that penalties which only depend on the number of parameters of each model, as for the classical criteria, are theoretically (and also practically) not adapted. This was suggested by [16, 7] who show that the penalty term needs to be well defined, and in particular needs to depend on the complexity of the list of models, i.e. the number of models having the same dimension. For this reason, following the work of [10] and in particular [17] in the density estimation framework, we consider a penalized log-likelihood procedure to estimate the true distribution s of a Poisson or negative binomial-distributed sequence \mathbf{y} . We prove that, up to a $\log n$ factor, the resulting estimator satisfies an oracle inequality.

The paper is organized as follows. The general framework is described in Section

1. More precisely, we present our proposed penalized maximum-likelihood estimator, the form of the penalty and give some non-asymptotic risk bounds for the resulting estimator. The studies of the two considered models (Poisson and negative binomial) are done in parallel along the paper. Some exponential bounds are derived in Section 2. A simulation study is performed to compare our proposed criterion with others and an application to the segmentation of RNA-seq data illustrates the procedure in Section 3. The proof of the main result is given in Section 4 for which the proofs of some intermediate results are given in the Appendix 5.

1 Model Selection Procedure

1.1 Penalized maximum-likelihood estimator

Let us denote by m a partition of $\llbracket 1, n \rrbracket$, $m = \{\llbracket 1, \tau_1 \rrbracket, \llbracket \tau_1, \tau_2 \rrbracket, \dots, \llbracket \tau_k, n \rrbracket\}$ and by \mathcal{M}_n a set of partitions of $\llbracket 1, n \rrbracket$. In our framework we want to estimate the distribution s defined by $s(t) = \mathcal{G}(\theta_t, \phi)$, $1 \leq t \leq n$, and we consider the two following models:

$$\begin{aligned} \mathcal{G}(\theta_t, \phi) &= \mathcal{P}(\lambda_t) & (\mathcal{P}) \\ \mathcal{G}(\theta_t, \phi) &= \mathcal{NB}(p_t, \phi) & (\mathcal{NB}) \end{aligned}$$

In the (\mathcal{NB}) case, we suppose that the over-dispersion parameter ϕ is known. We define the collection of models :

Definition 1.1. The collection of models associated to partition m is \mathcal{S}_m the set of distribution of sequences of length n such that for each element s_m of \mathcal{S}_m , for each segment J of m , and for each t in J , $s_m(t) = \mathcal{G}(\theta_J, \phi)$:

$$\mathcal{S}_m = \{s_m \mid \forall J \in m, \forall t \in J, s_m(t) = \mathcal{G}(\theta_J, \phi)\}.$$

We shall denote by $|m|$ the number of segments in partition m , and by $|J|$ the length of segment J .

We consider the log-likelihood contrast $\gamma(u) = \sum_{t=1}^n -\log \mathbf{P}_u(Y_t)$, namely respectively for $u(t) = \mathcal{P}(\mu_t)$ and $u(t) = \mathcal{NB}(q_t, \phi)$,

$$\begin{aligned} \gamma(u) &= \sum_{t=1}^n \mu_t - Y_t \log(\mu_t) + \log(Y_t!), & (\mathcal{P}) \\ \gamma(u) &= \sum_{t=1}^n -\phi \log q_t - Y_t \log(1 - q_t) - \log \left(\frac{\Gamma(\phi + Y_t)}{\Gamma(\phi) Y_t!} \right). & (\mathcal{NB}) \end{aligned}$$

Then the minimal contrast estimator \hat{s}_m of s on the collection \mathcal{S}_m is

$$\hat{s}_m = \arg \min_{u \in \mathcal{S}_m} \gamma(u), \tag{1}$$

so that, noting $\bar{Y}_J = \frac{\sum_{t \in J} Y_t}{|J|}$, for all $J \in m$ and $t \in J$

$$\hat{s}_m(t) = \mathcal{P}(\bar{Y}_J) \text{ for } (\mathcal{P}) \quad \text{and} \quad \hat{s}_m(t) = \mathcal{NB}\left(\frac{\phi}{\phi + \bar{Y}_J}, \phi\right) \text{ for } (\mathcal{NB}). \tag{2}$$

Therefore, for each partition m of \mathcal{M}_n we can obtain the best estimator \hat{s}_m as in equation (2), and thus define a collection of estimators $\{(\hat{s}_m)_{m \in \mathcal{M}_n}\}$. Ideally, we would wish to select the estimator $\hat{s}_{m(s)}$ amongst this collection with the minimum given risk. In the log-likelihood framework, it is natural to consider the Kullback-Leibler risk, with $K(s, u) = \mathbf{E}[\gamma(u) - \gamma(s)]$. In the following we note \mathbf{E} and \mathbf{P} the expectation and the probability under the true distribution s respectively (otherwise the underlying distribution is mentioned). In our models, the Kullback-Leibler between distributions s and u can be developed into

$$K(s, u) = \sum_{t=1}^n \left(\mu_t - \lambda_t - \lambda_t \log \frac{\mu_t}{\lambda_t} \right), \quad (\mathcal{P})$$

$$K(s, u) = \phi \sum_{t=1}^n \log \left(\frac{p_t}{q_t} \right) + \frac{1 - p_t}{p_t} \log \left(\frac{1 - p_t}{1 - q_t} \right). \quad (\mathcal{NB})$$

Unfortunately, minimizing this risk requires the knowledge of the true distribution s , and is unreachable. We will therefore want to consider the estimator $\hat{s}_{\hat{m}}$ where \hat{m} minimizes $\gamma(\hat{s}_m) + \text{pen}(m)$ for a well-chosen function pen (depending on the data). By doing so, we hope to select an estimator $\hat{s}_{\hat{m}}$ whose risk is as close as possible to the risk of $\hat{s}_{m(s)} = \arg \min_{m \in \mathcal{M}_n} \mathbf{E}_s[K(s, \hat{s}_m)]$ in the sense that

$$\mathbf{E}[K(s, \hat{s}_{\hat{m}})] \leq C \mathbf{E}[K(s, \hat{s}_{m(s)})],$$

where C is a nonnegative constant hopefully close to 1. We therefore introduce the following definition:

Definition 1.2. Let \mathcal{M}_n be a collection of partitions of $\llbracket 1, n \rrbracket$ constructed on a partition m_f (i.e. m_f is a refinement of every m in \mathcal{M}_n). Given a nonnegative, increasing in the size of m penalty function $\text{pen}: \mathcal{M}_n \rightarrow \mathbf{R}_+$, and choosing

$$\hat{m} = \arg \min_{m \in \mathcal{M}_n} \{\gamma(\hat{s}_m) + \text{pen}(m)\},$$

we define the penalized maximum-likelihood estimator as $\hat{s}_{\hat{m}}$.

In the following Section we provide a choice of penalty function, and show that the resulting estimator satisfies an oracle inequality.

1.2 Choice of the penalty function

Main result

The following result shows that for an appropriate choice of the penalty function, we have a non-asymptotic risk bound for the penalized maximum-likelihood estimator.

Theorem 1.3. *Let \mathcal{M}_n be a collection of partitions constructed on a partition m_f such that there exist absolute positive constants ρ_{\min} , ρ_{\max} and Γ satisfying:*

- $\forall t, \rho_{\min} \leq \theta_t \leq \rho_{\max}$ and

- $\forall J \in m_f, |J| \geq \Gamma(\log(n))^2$.

Let $(L_m)_{m \in \mathcal{M}_n}$ be some family of positive weights satisfying

$$\Sigma = \sum_{m \in \mathcal{M}_n} \exp(-L_m |m|) < +\infty. \quad (3)$$

Let $\beta > 1/2$ in the Poisson case, $\beta > 1/4$ in the negative binomial case. If for every $m \in \mathcal{M}_n$

$$\text{pen}(m) \geq \beta |m| \left(1 + 4\sqrt{L_m}\right)^2, \quad (4)$$

then

$$\mathbf{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C_\beta \inf_{m \in \mathcal{M}_n} \{K(s, \bar{s}_m) + \text{pen}(m)\} + C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, \Sigma),$$

with $C_\beta = \frac{(16\beta)^{1/3}}{(2\beta)^{1/3} - 1}$ in model (\mathcal{P}) and $C_\beta = \frac{(4\beta)^{1/3}}{(4\beta)^{1/3} - 1}$ in model (\mathcal{NB}) .

We note $h^2(s, u)$ the squared Hellinger distance between distribution s and u and \bar{s}_m is the projection of s onto the collection \mathcal{S}_m according to the Kullback-Leibler distance. The proof of this Theorem is given in Section 4.

Denoting $\bar{s}_m = \arg \min_{u \in \mathcal{S}_m} K(s, u)$, we have for $J \in m$ and $t \in J$,

$$\begin{aligned} \bar{s}_m(t) &= \mathcal{P}(\bar{\lambda}_J) & \text{where } \bar{\lambda}_J &= \frac{\sum_{t \in J} \lambda_t}{|J|} & (\mathcal{P}) \\ \bar{s}_m(t) &= \mathcal{NB}(p_J, \phi) & \text{where } p_J &= \frac{|J|}{\sum_{t \in J} 1/p_t}. & (\mathcal{NB}) \end{aligned} \quad (5)$$

We remark that the risk of the penalized estimator $\hat{s}_{\hat{m}}$ is treated in terms of Hellinger distance instead of the Kullback-Leibler information. This is due to the fact that the Kullback-Leibler is possibly infinite, and so difficult to control. It is possible to obtain a risk bound in term of Kullback-Leibler if we have a uniform control of $\|\log(s/\bar{s}_m)\|_\infty$ (see [18] for more explanation).

Choice of the weights $\{L_m, m \in \mathcal{M}_n\}$.

The penalty function depends on the family \mathcal{M}_n through the choice of the weights L_m which satisfy (3). We consider for \mathcal{M}_n the set of all possible partitions of $\llbracket 1, n \rrbracket$ constructed on a partition m_f which satisfies, for all segment J in m_f , $|J| \geq \Gamma(\log n)^2$. Classically (see [19]) the weights are chosen as a function of the dimension of the model s , which is here $|m|$. The number of partitions of \mathcal{M}_n having dimension D being bounded by $\binom{n}{D}$,

we have

$$\begin{aligned}
\Sigma &= \sum_{m \in \mathcal{M}_n} e^{L_m |m|} = \sum_{D=1}^n e^{-L_D D} \text{Card}\{m \in \mathcal{M}_n, |m| = D\} \\
&\leq \sum_{D=1}^n \binom{n}{D} e^{-L_D D} \leq \sum_{D=1}^n \left(\frac{en}{D}\right)^D e^{-L_D D} \\
&\leq \sum_{D=1}^n e^{-D \left(L_D - 1 - \log\left(\frac{n}{D}\right)\right)}.
\end{aligned}$$

So with the choice $L_D = 1 + \kappa + \log\left(\frac{n}{D}\right)$ with $\kappa > 0$, condition (3) is satisfied. Choosing, say $\kappa = 0.1$, the penalty function can be chosen of the form

$$\text{pen}(m) = \beta |m| \left(1 + 4 \sqrt{1.1 + \log\left(\frac{n}{|m|}\right)}\right)^2, \quad (6)$$

where β is a constant to be calibrated.

Integrating this penalty in Theorem 1.3 leads to the following control:

$$\mathbf{E} [h^2(s, \hat{s}_m)] \leq C_1 \inf_{m \in \mathcal{M}_n} \left\{ K(s, \bar{s}_m) + \beta |m| \left(1 + 4 \sqrt{1.1 + \log\left(\frac{n}{|m|}\right)}\right)^2 \right\} + C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, \Sigma) \quad (7)$$

The following proposition gives a bound on the Kullback-Leibler risk associated to \hat{s}_m :

Proposition 1.4. *Let m be a partition of \mathcal{M}_n , \hat{s}_m be the minimum contrast estimator and \bar{s}_m be the projection of s given by equations (2) and (5) respectively. Assume that there exists some positive absolute constants ρ_{\min} , ρ_{\max} and Γ such that $\forall t, \rho_{\min} \leq \theta_t \leq \rho_{\max}$ and $|J| \geq \Gamma(\log n)^2$. Then $\forall \varepsilon > 0, \forall a > 2$*

$$K(s, \bar{s}_m) - \frac{C_1(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \varepsilon, a)}{n^{a/2-\alpha}} + C_2(\varepsilon)|m| \leq \mathbf{E}[K(s, \hat{s}_m)],$$

where $\alpha < 1$ is a constant that can be expressed according to n , $C_2(\varepsilon) = \frac{1}{2} \frac{1-\varepsilon}{(1+\varepsilon)^2}$ in the

Poisson model (\mathcal{P}) and $C_2(\varepsilon) = \rho_{\min}^2 \frac{(1-\varepsilon)^2}{(1+\varepsilon)^4}$ in the negative binomial model (\mathcal{NB}).

The proof is given in appendix 5.1.

Combining proposition 1.4 and equation (7), we obtain the following oracle-type inequality:

Corollary 1.5. *Let \mathcal{M}_n be a collection of partitions constructed on a partition m_f such that there exist absolute positive constants ρ_{\min} , ρ_{\max} and Γ verifying:*

- $\forall t, \rho_{\min} \leq \theta_t \leq \rho_{\max}$ and
- $\forall J \in \mathcal{M}_f, |J| \geq \Gamma(\log n)^2$.

There exists some absolute constant C such that

$$\mathbf{E}[h^2(s, \hat{s}_m)] \leq C \log(n) \inf_{m \in \mathcal{M}_n} \{\mathbf{E}[K(s, \hat{s}_m)]\} + C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, \Sigma).$$

2 Exponential bounds

In order to prove Theorem 1.3, the general procedure in this model selection framework (see for example [19]) is the following: by definitions of \hat{m} and \hat{s}_m (see definition 1.2 and equation (1)), we have, $\forall m \in \mathcal{M}_n$

$$\gamma(\hat{s}_m) + \text{pen}(\hat{m}) \leq \gamma(\hat{s}_m) + \text{pen}(m) \leq \gamma(\bar{s}_m) + \text{pen}(m).$$

Then, with $\bar{\gamma}(u) = \gamma(u) - \mathbf{E}[\gamma(u)]$,

$$K(s, \hat{s}_m) \leq K(s, \bar{s}_m) + \bar{\gamma}(\bar{s}_m) - \bar{\gamma}(\hat{s}_m) - \text{pen}(\hat{m}) + \text{pen}(m).$$

The idea is therefore to control $\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(\hat{s}_{m'})$ uniformly over $m' \in \mathcal{M}_n$. This is more complicated when dealing with different models m and m' . Thus, following the work of [17] (see proof of Theorem 3.2, also recalled in [18]), we propose the following decomposition

$$\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(\hat{s}_{m'}) = (\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'})) + (\bar{\gamma}(s) - \bar{\gamma}(\bar{s}_{m'})) + (\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s)), \quad (8)$$

and control each term separately. The first term is the most delicate to handle, and requires the introduction and the control of a chi-square statistic. The main difficulty here is the non-bounded characteristic of the objects we are dealing with. Indeed, in the classic density estimation context such as that of [17], the objects are probabilities which are bounded and so facilitate the direct use of concentration inequalities.

In our case, the chi-square statistic we introduce is denoted χ_m^2 and defined by

$$\chi_m^2 = \chi^2(\bar{s}_m, \hat{s}_m) = \sum_{J \in m} |J| \frac{(\bar{Y}_J - \bar{E}_J)^2}{\bar{E}_J}, \quad (9)$$

where we recall that $\bar{Y}_J = \frac{\sum_{t \in J} Y_t}{|J|}$ and use the notation $\bar{E}_J = \frac{E_J}{|J|}$ with $E_J = \sum_{t \in J} E_t$.

Respectively for (\mathcal{P}) and (\mathcal{NB}) , we have $E_t = \lambda_t$ and $E_t = \phi \frac{1-p_t}{p_t}$. The purpose is thus to control χ_m^2 uniformly over \mathcal{M}_n . To this effect, we need to obtain an exponential bound of $Y_J = \sum_{t \in J} Y_t$ around its expectation. In Subsection 2.1, we recall a result of [15] that we use to derive an exponential bound for χ_m^2 (Subsection 2.2).

2.1 Control of Y_J

First we recall a large deviation results established by [15] (lemma 3) that we apply in the Poisson and negative binomial frameworks.

Lemma 2.1. *Let Y_1, \dots, Y_n be n independent centered random variables.*

If $\log(\mathbf{E}[e^{zY_i}]) \leq \kappa \frac{z^2 \theta_i}{2(1-z\tau)}$ for all $z \in [0, 1/\tau[$, and $1 \leq i \leq n$, then

$$\mathbf{P} \left[\sum_{i=1}^n Y_i \geq \left(2\kappa x \sum_{i=1}^n \theta_i \right)^{1/2} + \tau x \right] \leq e^{-x} \text{ for all } x > 0.$$

If for $1 \leq i \leq n$ and all $z > 0$ $\log(\mathbf{E}[e^{-zY_i}]) \leq \kappa z^2 \theta_i / 2$, then

$$\mathbf{P} \left[\sum_{i=1}^n Y_i \leq - \left(2\kappa x \sum_{i=1}^n \theta_i \right)^{1/2} \right] \leq e^{-x} \text{ for all } x > 0.$$

To apply this lemma we therefore need a majoration of $\log \mathbf{E} [e^{z(Y_t - E_t)}]$ and $\log \mathbf{E} [e^{-z(Y_t - E_t)}]$ for $z > 0$.

Poisson case.

With $E_t = \lambda_t$, we have:

$$\log \mathbf{E} [e^{z(Y_t - \lambda_t)}] = -z\lambda_t + \log \mathbf{E} [e^{zY_J}] = -z\lambda_t + \log e^{(\lambda_t(e^z - 1))} = \lambda_t(e^z - z - 1).$$

So

$$\log \mathbf{E} [e^{z(Y_t - E_t)}] = E_t(e^z - z - 1).$$

Negative binomial case.

In this case $E_t = \phi \frac{1-p_t}{p_t}$ and we have

$$\begin{aligned} \log \mathbf{E} \left(e^{z(Y_t - \phi \frac{1-p_t}{p_t})} \right) &= -z\phi \frac{1-p_t}{p_t} + \phi \log \frac{p_t}{1 - (1-p_t)e^z} \text{ for } z \leq -\log(1-p_t) \\ &\leq \phi \left[\frac{1-p_t}{p_t}(-z) + \frac{1-p_t}{p_t} \frac{p_t}{1 - (1-p_t)e^z} (e^z - 1) \right] \\ &\leq \phi \left[\frac{1-p_t}{p_t}(-z) + \frac{1-p_t}{p_t} (e^z - 1) \right] \leq \phi \frac{1-p_t}{p_t} (e^z - z - 1). \end{aligned}$$

So that in both cases,

$$\log \mathbf{E} [e^{z(Y_t - E_t)}] \leq E_t(e^z - z - 1).$$

Now using $e^z - z - 1 \leq \frac{z^2}{2(1-z)}$ for $z > 0$ and $e^z - z - 1 \leq \frac{z^2}{2}$ for $z < 0$, we have

$$\log \mathbf{E} [e^{z(Y_t - E_t)}] \leq E_t \frac{z^2}{2(1-z)} \quad \text{and} \quad \log \mathbf{E} [e^{-z(Y_t - E_t)}] \leq E_t \frac{z^2}{2}$$

Then,

$$P [Y_J - E_J \geq \sqrt{2xE_J} + x] \leq e^{-x},$$

or

$$P [Y_J - E_J \geq x] \leq e^{-\frac{x^2}{2(E_J+x)}} \quad \text{and} \quad P [|Y_J - E_J| \geq x] \leq 2e^{-\frac{x^2}{2(E_J+x)}} \quad (10)$$

2.2 Exponential bound for χ_m^2

We first introduce the following set Ω_m defined by:

$$\Omega_m(\varepsilon) = \bigcap_{J \in m} \left\{ \left| \frac{Y_J}{E_J} - 1 \right| \leq \varepsilon \right\}, \quad (11)$$

for all $\varepsilon \in]0, 1[$ and all segmentations m such that each segment J verifies $|J| \geq \Gamma(\log(n))^2$. This set has a large probability since we obtain

$$\begin{aligned} \mathbf{P}(\Omega_m(\varepsilon)^C) &\leq \sum_{J \in m} \mathbf{P}(|Y_J - E_J| > \varepsilon E_J) \leq 2 \sum_{J \in m} e^{-\frac{\varepsilon^2 E_J}{2(1+\varepsilon)}} \\ &\leq 2 \sum_{J \in m} e^{-|J| \varepsilon' f(\phi, \rho_{\min})} \leq 2|m| \exp(-\varepsilon' \Gamma f(\phi, \rho_{\min}) (\log(n))^2) \end{aligned}$$

by applying equation (10) with $x = \varepsilon E_J$ and where $\varepsilon' = \varepsilon^2 / (2(1+\varepsilon))$ and $f(\phi, \rho_{\min}) > 0$. Thus

$$\mathbf{P}(\Omega_m(\varepsilon)^C) \leq \frac{C(\phi, \Gamma, \rho_{\min}, \varepsilon, a)}{n^a}, \quad (12)$$

with $a > 2$.

The reason for introducing this set is double: in addition to enable the control of χ_m^2 given by equation (9) on this restricted set, it allows us to link $K(\hat{s}_m, \bar{s}_m)$ to V_m^2 (see (18) for the control of the first term in the decomposition) and so to χ_m^2 , relation that we use to evaluate the risk of one model (see (20)).

Let m_f be a partition of \mathcal{M}_n such that $\forall J \in m_f, |J| \geq \Gamma(\log(n))^2$ and assume that all considered partitions in \mathcal{M}_n are constructed on this grid m_f . The following proposition gives an exponential bound for χ_m^2 on the restricted event $\Omega_{m_f}(\varepsilon)$.

Proposition 2.2. *Let Y_1, \dots, Y_n be independent random variables with distribution \mathcal{G} (Poisson or negative binomial distribution). Let m be a partition of \mathcal{M}_n with $|m|$ segments and χ_m^2 the statistic given by (9). For any positive x , we have*

$$\mathbf{P} \left[\chi_m^2 \mathbf{1}_{\Omega_{m_f}(\varepsilon)} \geq |m| + 8(1+\varepsilon)\sqrt{x|m|} + 4(1+\varepsilon)x \right] \leq e^{-x}.$$

Proof. As in the density estimation framework, this quantity can be controlled using the Bernstein inequality. In our context, noting $\chi_m^2 = \sum_{J \in m} Z_J$ where

$$Z_J = \frac{(Y_J - E_J)^2}{E_J},$$

we need

- the calculation (or bounds) of the expectation of χ_m^2 :

Poisson case

Y_J is distributed according to a Poisson distribution with parameter λ_J so that

$$\mathbf{E} [\chi_m^2] = |m|. \quad (13)$$

Negative binomial case

We have

$$\mathbf{E} [\chi_m^2] = \sum_{J \in m} \frac{1}{|J|} \frac{\sum_{t \in J} \text{Var}(Y_t)}{\phi^{\frac{1-p_J}{p_J}}} = \sum_{J \in m} \frac{1}{|J|} \frac{\sum_{t \in J} \phi^{\frac{1-p_t}{p_t^2}}}{\phi^{\frac{1-p_J}{p_J}}},$$

and thus

$$|m| \leq \mathbf{E} [\chi_m^2] \leq \frac{1}{\rho_{\min}} |m|. \quad (14)$$

- an upper bound of $\sum_{J \in m} \mathbf{E}[Z_J^p]$. For every $p \geq 2$ we have,

$$\begin{aligned} \mathbf{E} [Z_J^p \mathbf{1}_{\Omega_{m_f}(\epsilon)}] &= \frac{1}{E_J^p} \int_0^{+\infty} 2p x^{2p-1} P[\{|Y_J - E_J| \geq x\} \cap \Omega_{m_f}(\epsilon)] dx \\ &\leq \frac{1}{E_J^p} \int_0^{\varepsilon E_J} 2p x^{2p-1} P[|Y_J - E_J| \geq x] dx. \end{aligned}$$

Using equation (10) and since $x \leq \varepsilon E_J$, we obtain the exponential bound $P[|Y_J - E_J| \geq x] \leq 2e^{-\frac{x^2}{2E_J(1+\varepsilon)}}$.

Therefore

$$\begin{aligned} \mathbf{E} [Z_J^p \mathbf{1}_{\Omega_{m_f}(\epsilon)}] &\leq \frac{1}{E_J^p} \int_0^{\varepsilon E_J} 4p x^{2p-1} e^{-\frac{x^2}{2E_J(1+\varepsilon)}} dx \\ &\leq 4p (1+\varepsilon)^p \int_0^{+\infty} u^{2p-1} e^{-\frac{u^2}{2}} du \\ &\leq 4p (1+\varepsilon)^p \int_0^{+\infty} (2t)^{p-1} e^{-t} dt \\ &\leq 2^{p+1} p (1+\varepsilon)^p p!, \end{aligned}$$

and

$$\sum_{J \in m} \mathbf{E} [Z_J^p \mathbf{1}_{\Omega_{m_f}(\epsilon)}] \leq 2^{p+1} p (1+\varepsilon)^p p! |m|.$$

Since $p \leq 2^{p-1}$,

$$\sum_{J \in m} \mathbf{E} [Z_J^p \mathbf{1}_{\Omega_{m_f}(\epsilon)}] \leq \frac{p!}{2} \times [2^5 (1+\varepsilon)^2 |m|] \times [4(1+\varepsilon)]^{p-2}.$$

We conclude by taking $v = 2^5 (1 + \varepsilon)^2 |m|$ and $c = 4(1 + \varepsilon)$ (see proposition 2.9 of [18] for the definition of the Bernstein's inequality).

□

3 Simulations and application

In the context of RNA-seq experiments, an important question is the (re)-annotation of the genome, that is, the precise localisation of the transcribed regions on the chromosomes. In an ideal situation, when considering the number of reads starting at each position, one would expect to observe a uniform coverage over each gene (proportional to its expression level), separated by regions of null signal (corresponding to non-transcribed regions of the genome). In practice however, those experiments tend to return very noisy signals that are best modelled by the negative binomial distribution.

In this Section, we first study the performance of the proposed penalized criterion by comparing it with others model selection criteria on a resampling dataset (Subsection 3.1). Then we provide an application on real data (Subsection 3.2). Since the penalty depends on the partition only through its size, the segmentation procedure is two-steps: first we estimate, for all number of segments K between 1 and K_{max} , the optimal partition with K segments (i.e. construct the collection of estimators $\{\hat{s}_K\}_{1 \leq K \leq K_{max}}$ where $\hat{s}_K = \arg \min_{\hat{s}_m, m \in \mathcal{M}_K} \{\gamma(\hat{s}_m)\}$). The optimal solution is obtained using a fast segmentation algorithm such as the Pruned Dynamic Programming Algorithm (PDPA, [20]) implemented for the Poisson and negative binomial losses or contrasts in the R package `Segmentor3IsBack` [21]. Then, we choose K using our penalty function which requires the calibration of the constant β that can be tuned according to the data by using the slope heuristic (see [7, 22]). Using the negative binomial distribution requires the knowledge of parameter ϕ . We propose to estimate it using a modified version of the Jonhson and Kotz's estimator [23].

3.1 Simulation study

We have assessed the performances of the proposed method (called Penalized PDPA) on a simulation scenario by comparing to five other procedures both its choice in the number of segments and the quality of the obtained segmentation using the Rand-Index \mathcal{I} . This index is defined as follows: let C_t be the true index of the segment to which base t belongs and let \hat{C}_t be the corresponding estimated index, then

$$\mathcal{I} = \frac{2 \sum_{t>s} \left[\mathbf{1}_{C_t=C_s} \mathbf{1}_{\hat{C}_t=\hat{C}_s} + \mathbf{1}_{C_t \neq C_s} \mathbf{1}_{\hat{C}_t \neq \hat{C}_s} \right]}{(n-1)(n-2)}.$$

The characteristics of the different algorithms are described in Table 1.

The data we considered comes from a resampling procedure using real RNA-seq data. The original data, from a study by the Sherlock Genomics laboratory at Stanford University, is publicly available on the NCBI's Sequence Read Archive (SRA, url: <http://www.ncbi.nlm.nih.gov/sra>) with the accession number SRA048710. We created

Algorithm	Dist	Complexity	Inference	Pen	Exact	Reference
Penalized PDPA	NB	$n \log n$	frequentist	external	exact	[21]
PDPA with BIC	NB	$n \log n$	frequentist	external	exact	[21]
Penalized PDPA	P	$n \log n$	frequentist	external	exact	[21]
PDPA with BIC	P	$n \log n$	frequentist	external	exact	[21]
PELT with BIC	P	n	frequentist	internal	exact	[24]
CART with BIC	P	$n \log n$	frequentist	external	heuristic	[25]
postCP with ICL	NB	n	frequentist	external	exact	[4]
EBS with ICL	NB	n^2	Bayesian	external	exact	[26]

Table 1: *Properties of segmentation algorithms.* The first column indicates the name of the algorithm and the criterion used for the choice of K . In the second column, NB stands for the negative binomial distribution and P for Poisson. The time of each algorithm is given (column "Complexity") and column "Exact" precises if the exact solution is reached.

an artificial gene, inspired from the *Drosophila* *inr-a* gene, resulting in a 14-segment signal with unregular intensities mimicking a differentially transcribed gene. 100 datasets are thus created. Results are presented using boxplots in Figure 3.1. Because PELT's estimate of K averaged around 427 segments, we did not show its corresponding boxplot.

We can see that with the negative binomial distribution, not only do we perfectly recover the true number of segments, but our procedure outperforms all other approaches. Moreover, the impressive results in terms of Rand-Index prove that our choice of number of segments also leads to the almost perfect recovery of the true segmentation. However, the use of the Poisson loss leads to a constant underestimation of the number of segments, which is reflected on the Rand-Index values. This is due to the inappropriate choice of distribution (confirmed by the other algorithms implemented for the Poisson loss which perform worse than the others). It however underlines the need for the development of methods for the negative binomial distribution. Moreover, in terms of computational time, the fast algorithm [21] is in $\mathcal{O}(n \log n)$, allowing its use on long signals (such as a whole-genome analysis), even though it is not as fast as CART or PELT.

3.2 Segmentation of RNA-Seq data

We apply our proposed procedure for segmenting chromosome 1 of the *S. Cerevisiae* (yeast) using RNA-Seq data from the Sherlock Laboratory at Stanford University [1] and publicly available from the NCBI's Sequence Read Archive (SRA, url:<http://www.ncbi.nlm.nih.gov/sra>, accession number SRA048710). An existing annotation is available on the *Saccharomyces* Genome Database (SGD) at url:<http://www.yeastgenome.org>, which allows us to validate our results. The two distributions (Poisson and negative binomial) are considered here to show the difference.

In the Poisson distribution case, we select 106 segments of which only 19 are related to the SGD annotation. Indeed, as illustrated by Figure 2, 36 of the segments have a length smaller than 10: the Poisson loss is not adapted to this kind of data with high variability and it tends to select outliers as segment. On the contrary, we select 103 segments in the

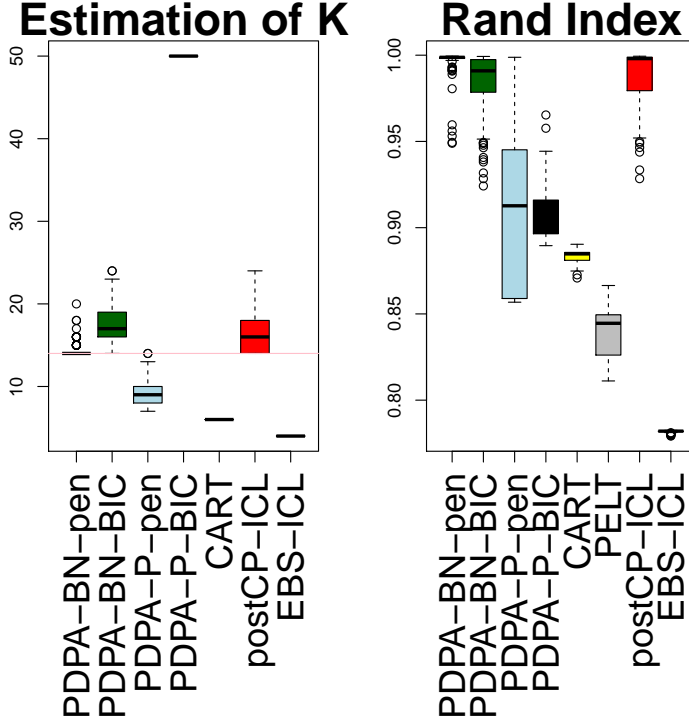


Figure 1: **Estimation of K on resampled datasets.** Left: boxplot of the estimation of K on data-sets simulated by resampling on artificial gene Inr-a. PELT’s estimates average at 427 segments and are not shown. The pink horizontal line indicates the true value of K . Right: boxplot of the Rand-Index values for the proposed estimators.

negative binomial case most of which (all but 3) surround known genes from the SGD. Figure 3 illustrates the result. However, almost none of those change-points correspond exactly to annotated boundaries. Discussion with biologists has increased our belief in the need for genome (re-)annotation using RNA-seq data, and in the validity of our approach.

4 Proof of Theorem 1.3

Recall that we want to control the three terms in the decomposition given by (8). All the proofs of the different propositions are given in Section 5.

- The control the term $\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'})$ is obtained with the following proposition where the set $\Omega_1(\xi)$ is defined by

$$\Omega_1(\xi) = \bigcap_{m' \in \mathcal{M}_n} \left\{ \chi_{m'}^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq |m'| + 8(1 + \epsilon) \sqrt{(L_{m'} |m'| + \xi) |m'|} + 4(1 + \epsilon)(L_{m'} |m'| + \xi) \right\}.$$

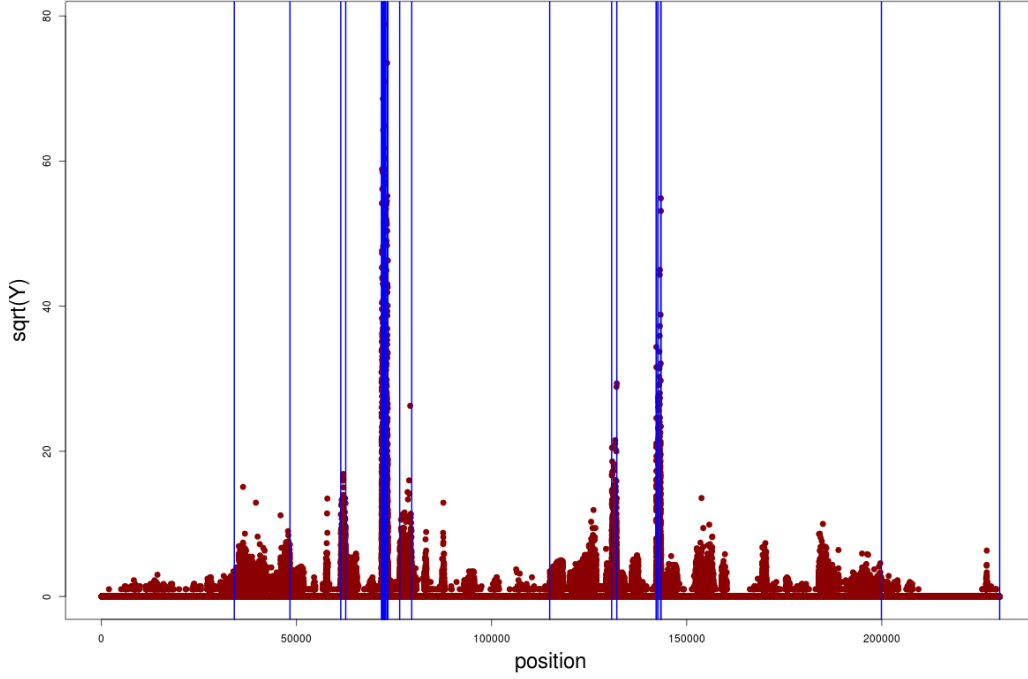


Figure 2: *Segmentation of the yeast chromosome 1 using Poisson loss.* Read-count are represented on a root-squared scale. The model selection procedure chooses $K = 106$ segments.

Proposition 4.1. *Let m' be a partition of \mathcal{M}_n . Then*

$$\begin{aligned} (\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon) \cap \Omega_1(\xi)} &\leq C(\epsilon) \left[|m'| + 8(1 + \epsilon) \sqrt{(L_{m'}|m'| + \xi)|m'|} \right. \\ &\quad \left. + 4(1 + \epsilon)(L_{m'}|m'| + \xi) \right] + \frac{1}{1 + \epsilon} K(\bar{s}_{m'}, \hat{s}_{m'}), \end{aligned}$$

with $C(\epsilon) = \frac{1}{2} \left(\frac{1+\epsilon}{1-\epsilon} \right)$ in the Poisson case and $C(\epsilon) = \frac{1+\epsilon}{4}$ in the negative binomial case.

- The control of the term $\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s)$, or more precisely its expectation, is given by the following proposition:

Proposition 4.2.

$$|\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s)) \mathbf{1}_{\Omega_{m_f}(\epsilon)}]| \leq \frac{C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \epsilon, a)}{n^{(a-1)/2}}. \quad (15)$$

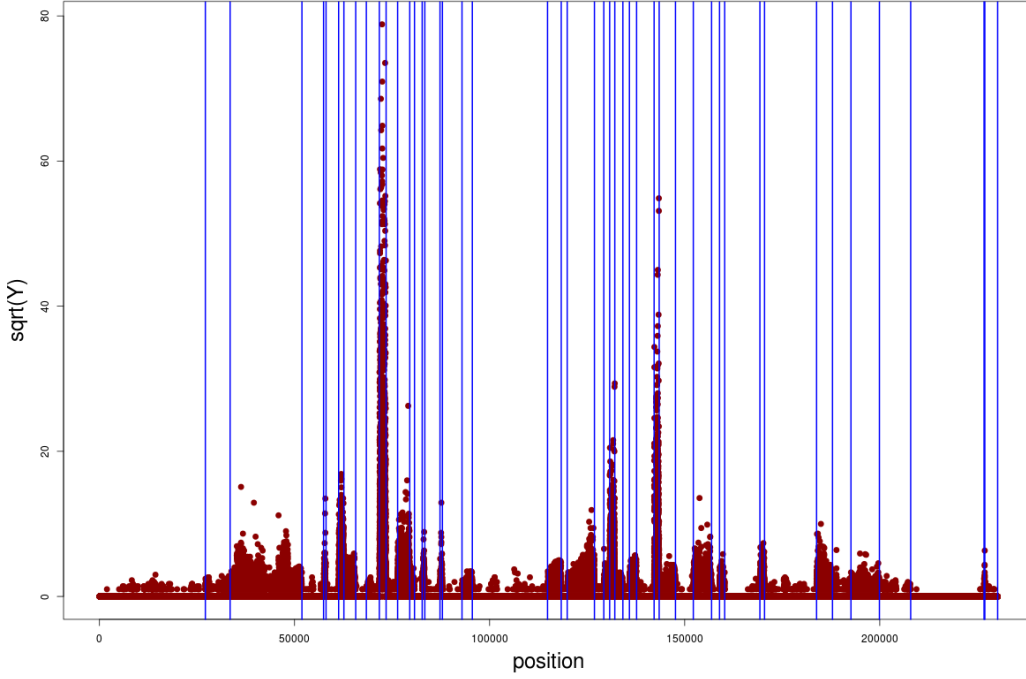


Figure 3: *Segmentation of the yeast chromosome 1 using the negative binomial loss.* The model selection procedure chooses $K = 103$ segments, most of which surround genes given by the SGD annotation.

- To control $\bar{\gamma}(s) - \bar{\gamma}(\bar{s}_{m'})$, we use the following proposition which gives an exponential bound for $\bar{\gamma}(s) - \bar{\gamma}(u)$.

Proposition 4.3. *Let s and u be two distributions of a sequence Y . Let γ be the log-likelihood contrast, $\bar{\gamma}(u) = \gamma(u) - \mathbf{E}[\gamma(u)]$, and $K(s, u)$ and $h^2(s, u)$ be respectively the Kullback-Leibler and the squared Hellinger distances between distributions s and u . Then $\forall x > 0$,*

$$\mathbf{P} [\bar{\gamma}(s) - \bar{\gamma}(u) \geq K(s, u) - 2h^2(s, u) + 2x] \leq e^{-x}.$$

Applying it to $u = \bar{s}_{m'}$ yields:

$$\mathbf{P} [\bar{\gamma}(s) - \bar{\gamma}(\bar{s}_{m'}) \geq K(s, \bar{s}_{m'}) - 2h^2(s, \bar{s}_{m'}) + 2x] \leq e^{-x}. \quad (16)$$

We then define $\Omega_2(\xi) = \bigcap_{m' \in \mathcal{M}_n} \{\bar{\gamma}(s) - \bar{\gamma}(\bar{s}_{m'}) \leq K(s, \bar{s}_{m'}) - 2h^2(s, \bar{s}_{m'}) + 2(L_{m'}|m'| + \xi)\}$.

Let $\Omega(\varepsilon, \xi) = \Omega_{m_f}(\varepsilon) \cap \Omega_1(\xi) \cap \Omega_2(\xi)$. Then, combining equation (16) and proposition 4.1, we get for $m' = \hat{m}$,

$$\begin{aligned}
(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(\hat{s}_{\hat{m}}))\mathbf{1}_{\Omega(\epsilon, \xi)} &= (\bar{\gamma}(s) - \bar{\gamma}(\bar{s}_{\hat{m}}))\mathbf{1}_{\Omega(\epsilon, \xi)} + (\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega(\epsilon, \xi)} + (\bar{\gamma}(\bar{s}_{\hat{m}}) - \bar{\gamma}(\hat{s}_{\hat{m}}))\mathbf{1}_{\Omega(\epsilon, \xi)} \\
&\leq [K(s, \bar{s}_{\hat{m}}) - 2h^2(s, \bar{s}_{\hat{m}})]\mathbf{1}_{\Omega(\epsilon, \xi)} + R\mathbf{1}_{\Omega(\epsilon, \xi)} + \frac{1}{1+\varepsilon}K(\bar{s}_{\hat{m}}, \hat{s}_{\hat{m}})\mathbf{1}_{\Omega(\epsilon, \xi)} \\
&\quad + C(\varepsilon) \left[|\hat{m}| + 8(1+\varepsilon)\sqrt{(L_{\hat{m}}|\hat{m}| + \xi)|\hat{m}|} + 4(1+\varepsilon)(L_{\hat{m}}|\hat{m}| + \xi) \right] \\
&\quad + 2L_{\hat{m}}|\hat{m}| + 2\xi,
\end{aligned}$$

with $R = \bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s)$. So that

$$\begin{aligned}
K(s, \hat{s}_{\hat{m}})\mathbf{1}_{\Omega(\epsilon, \xi)} &\leq [K(s, \bar{s}_{\hat{m}}) - 2h^2(s, \bar{s}_{\hat{m}})]\mathbf{1}_{\Omega(\epsilon, \xi)} + \frac{1}{1+\varepsilon}K(\bar{s}_{\hat{m}}, \hat{s}_{\hat{m}})\mathbf{1}_{\Omega(\epsilon, \xi)} \\
&\quad + C(\varepsilon) \left[|\hat{m}| + 8(1+\varepsilon)\sqrt{(L_{\hat{m}}|\hat{m}| + \xi)|\hat{m}|} + 4(1+\varepsilon)(L_{\hat{m}}|\hat{m}| + \xi) \right] \\
&\quad + K(s, \bar{s}_m)\mathbf{1}_{\Omega(\epsilon, \xi)} + 2L_{\hat{m}}|\hat{m}| + 2\xi + R\mathbf{1}_{\Omega(\epsilon, \xi)} - \text{pen}(\hat{m}) + \text{pen}(m).
\end{aligned}$$

And since

- $K(s, \hat{s}_{\hat{m}}) = K(s, \bar{s}_{\hat{m}}) + K(\bar{s}_{\hat{m}}, \hat{s}_{\hat{m}})$ (see equation (17)),
- $K(s, u) \geq 2h^2(s, u)$ (see lemma 7.23 in [18]),
- $h^2(s, \hat{s}_{\hat{m}}) \leq 2(h^2(s, \bar{s}_{\hat{m}}) + h^2(\bar{s}_{\hat{m}}, \hat{s}_{\hat{m}}))$ (using inequality $2ab \leq \kappa a^2 + \kappa^{-1}b^2$ with $\kappa = 1$),

$$\begin{aligned}
\frac{\varepsilon}{1+\varepsilon}h^2(s, \hat{s}_{\hat{m}})\mathbf{1}_{\Omega(\epsilon, \xi)} &\leq K(s, \bar{s}_m)\mathbf{1}_{\Omega(\epsilon, \xi)} + R\mathbf{1}_{\Omega(\epsilon, \xi)} - \text{pen}(\hat{m}) + \text{pen}(m) \\
&\quad + |\hat{m}|C(\varepsilon) \left[1 + (1+\varepsilon) \left(8\sqrt{L_{\hat{m}}} + \varepsilon + 4L_{\hat{m}} \right) \right] + 2L_{\hat{m}}|\hat{m}| \\
&\quad + 2\xi \left[1 + C(\varepsilon) \left(8(1+\varepsilon)\frac{2}{\varepsilon} + 4(1+\varepsilon) \right) \right].
\end{aligned}$$

But

$$\begin{aligned}
C(\varepsilon) \left[1 + (1+\varepsilon) \left(8\sqrt{L_{\hat{m}}} + \varepsilon + 4L_{\hat{m}} \right) \right] + 2L_{\hat{m}} &\leq C(\varepsilon) \left[1 + (1+\varepsilon) \left(\varepsilon + 8\sqrt{L_{\hat{m}}} + 8L_{\hat{m}} \right) \right] \\
&\leq C_2(\varepsilon) \left[1 + 8\sqrt{L_{\hat{m}}} + 8L_{\hat{m}} \right].
\end{aligned}$$

with $C_2(\varepsilon) = \frac{1}{2} \left(\frac{1+\varepsilon}{1-\varepsilon} \right)^3$ for (\mathcal{P}) and $C_2(\varepsilon) = \frac{1}{4} (1+\varepsilon)^3$ for (\mathcal{NB}) . So we have

$$\begin{aligned}
\frac{\varepsilon}{1+\varepsilon}h^2(s, \hat{s}_{\hat{m}})\mathbf{1}_{\Omega(\epsilon, \xi)} &\leq K(s, \bar{s}_m)\mathbf{1}_{\Omega(\epsilon, \xi)} + R\mathbf{1}_{\Omega(\epsilon, \xi)} - \text{pen}(\hat{m}) + \text{pen}(m) \\
&\quad + |\hat{m}|C_2(\varepsilon) \left(1 + 4\sqrt{L_{\hat{m}}} \right)^2 + 2\xi \left[1 + (1+\varepsilon)C(\varepsilon) \left(\frac{8}{\varepsilon} + 2 \right) \right].
\end{aligned}$$

By assumption, $\text{pen}(\hat{m}) \geq \beta|\hat{m}| (1 + 4\sqrt{L_{\hat{m}}})^2$. Choosing $\beta = C_2(\varepsilon)$ yields

$$h^2(s, \hat{s}_{\hat{m}}) \mathbf{1}_{\Omega(\epsilon, \xi)} \leq C_\beta [K(s, \bar{s}_m) \mathbf{1}_{\Omega(\epsilon, \xi)} + R \mathbf{1}_{\Omega(\epsilon, \xi)} + \text{pen}(m)] + \xi C(\beta).$$

Then, using propositions 4.2 and 4.1, we have $\mathbf{P}(\Omega_1(\xi)^C) \leq \sum_{m' \in \mathcal{M}_n} e^{-L_{m'}|m'|+\xi}$ and $\mathbf{P}(\Omega_2(\xi)^C) \leq \sum_{m' \in \mathcal{M}_n} e^{-L_{m'}|m'|+\xi}$. So that using hypothesis (3),

$$\mathbf{P}(\Omega_1(\xi)^C \cup \Omega_2(\xi)^C) \leq 2 \sum_{m' \in \mathcal{M}_n} e^{-L_{m'}|m'|+\xi} \leq 2\Sigma e^{-\xi},$$

and thus $\mathbf{P}(\Omega_1(\xi) \cap \Omega_2(\xi)) \geq 1 - 2\Sigma e^{-\xi}$. We now integrate over ξ , and using equation (15), we get with a probability larger than $1 - 2\Sigma e^{-\xi}$

$$\mathbf{E} \left[h^2(s, \hat{s}_{\hat{m}}) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \right] \leq C_\beta \left[K(s, \bar{s}_m) + \frac{C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, a)}{n^{(a-1)/2}} + \text{pen}(m) \right] + \Sigma C(\beta).$$

And since $\mathbf{E} \left[h^2(s, \hat{s}_{\hat{m}}) \mathbf{1}_{\Omega_{m_f}(\epsilon)^C} \right] \leq \frac{C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, a)}{n^{a-1}}$, we have

$$\mathbf{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C_\beta [K(s, \bar{s}_m) + \text{pen}(m)] + C'(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, \Sigma).$$

Finally, by minimizing over $m \in \mathcal{M}_n$, we get

$$\mathbf{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C_\beta \inf_{m \in \mathcal{M}_n} \{K(s, \bar{s}_m) + \text{pen}(m)\} + C'(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \beta, \Sigma).$$

5 Appendices

5.1 Proof of proposition 1.4

Using Pythagore-type identity, we obtain the following decomposition (see for example [17]):

$$K(s, \hat{s}_m) = K(s, \bar{s}_m) + K(\bar{s}_m, \hat{s}_m). \quad (17)$$

The objective is then to obtain a lower bound of $\mathbf{E}[K(\bar{s}_m, \hat{s}_m)]$ in the two considered distribution cases.

Poisson case

We have

$$K(\bar{s}_m, \hat{s}_m) = \sum_{J \in m} |J| \left(\bar{Y}_J - \bar{\lambda}_J - \bar{\lambda}_J \log \frac{\bar{Y}_J}{\bar{\lambda}_J} \right) = \sum_{J \in m} |J| \bar{\lambda}_J \Phi \left(\log \frac{\bar{Y}_J}{\bar{\lambda}_J} \right).$$

where $\Phi(x) = e^x - 1 - x$. Since $\frac{1}{2}x^2(1 \wedge e^x) \leq \Phi(x) \leq \frac{1}{2}x^2(1 \vee e^x)$, then on $\Omega_{m_f}(\epsilon)$, we have

$$\begin{aligned}\frac{1}{2} \log^2 \frac{\bar{Y}_J}{\bar{\lambda}_J} \left(1 \wedge \frac{\bar{Y}_J}{\bar{\lambda}_J}\right) &\leq \Phi \left(\log \frac{\bar{Y}_J}{\bar{\lambda}_J} \right) \leq \frac{1}{2} \log^2 \frac{\bar{Y}_J}{\bar{\lambda}_J} \left(1 \vee \frac{\bar{Y}_J}{\bar{\lambda}_J}\right), \\ \frac{1-\varepsilon}{2} \log^2 \frac{\bar{Y}_J}{\bar{\lambda}_J} &\leq \Phi \left(\log \frac{\bar{Y}_J}{\bar{\lambda}_J} \right) \leq \frac{1+\varepsilon}{2} \log^2 \frac{\bar{Y}_J}{\bar{\lambda}_J}.\end{aligned}$$

So

$$\frac{1-\varepsilon}{2} V_m^2 \leq K(\bar{s}_m, \hat{s}_m) \leq \frac{1+\varepsilon}{2} V_m^2, \quad (18)$$

where

$$V_m^2 = V^2(\bar{s}_m, \hat{s}_m) = \sum_{J \in m} |J| \bar{\lambda}_J \log^2 \frac{\bar{Y}_J}{\bar{\lambda}_J} = \sum_{J \in m} |J| \frac{(\bar{Y}_J - \bar{\lambda}_J)^2}{\bar{\lambda}_J} \left(\log \frac{\bar{Y}_J}{\bar{\lambda}_J} \right)^2. \quad (19)$$

And using, for $x > 0$, $\frac{1}{1 \vee x} \leq \frac{\log x}{x-1} \leq \frac{1}{1 \wedge x}$, we get, on $\Omega_{m_f}(\epsilon)$

$$\frac{1}{(1+\varepsilon)^2} \chi_m^2 \leq V_m^2 \leq \frac{1}{(1-\varepsilon)^2} \chi_m^2. \quad (20)$$

So

$$\frac{1-\varepsilon}{2(1+\varepsilon)^2} \chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq K(\bar{s}_m, \hat{s}_m) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq \frac{1+\varepsilon}{2(1-\varepsilon)^2} \chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)}.$$

On one hand, $\mathbf{E}[\chi_m^2] = |m|$, and

$$\frac{1-\varepsilon}{2(1+\varepsilon)^2} |m| - \mathbf{E} \left[\chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)^C} \right] \leq \mathbf{E} \left[K(\bar{s}_m, \hat{s}_m) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \right] \leq \frac{1+\varepsilon}{2(1-\varepsilon)^2} |m|.$$

Since $\chi_m^2 \leq \frac{1}{\Gamma(\log(n))^2 \rho_{\min}} \sum_{J \in m} (Y_J - \lambda_J)^2 \leq \frac{1}{\Gamma(\log(n))^2 \rho_{\min}} (\sum_t Y_t - \sum_t \lambda_t)^2$, using Cauchy-Schwarz Inequality, we get

$$\begin{aligned}\mathbf{E} \left[\chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)^C} \right] &\leq \frac{1}{\Gamma(\log(n))^2 \rho_{\min}} \left[3 \left(\sum_t \lambda_t \right)^2 + \sum_t \lambda_t \right]^{1/2} P(\Omega_{m_f}(\epsilon)^C)^{1/2} \\ &\leq C(\Gamma, \rho_{\min}, \rho_{\max}) \frac{n}{(\log(n))^2} P(\Omega_{m_f}(\epsilon)^C)^{1/2} \\ &\leq C(\Gamma, \rho_{\min}, \rho_{\max}) n^\alpha P(\Omega_{m_f}(\epsilon)^C)^{1/2} \\ &\leq \frac{C(\phi, \Gamma, \rho_{\min}, \rho_{\max}, \varepsilon, a)}{n^{a/2-\alpha}},\end{aligned}$$

where $\alpha = 1 - 2 \frac{\log(\log(n))}{\log(n)}$, $n \geq 2$. For example, $\alpha = 0.62$ for $n = 10^6$.

On the other hand, using $\log 1/x \geq 1 - x$ for all $x > 0$, $\mathbf{E} \left[K(\bar{s}_m, \hat{s}_m) \mathbf{1}_{\Omega_{m_f}(\epsilon)^C} \right] \geq 0$. Finally, we have

$$K(s, \bar{s}_m) + \frac{1-\varepsilon}{2(1+\varepsilon)^2} |m| - \frac{C_1(\Gamma, \rho_{\min}, \rho_{\max}, \varepsilon, a)}{n^{a/2-\alpha}} \leq \mathbf{E}[K(s, \hat{s}_m)],$$

Negative binomial case

We have $K(\bar{s}_m, \hat{s}_m) = \phi \sum_{J \in m} \frac{|J|}{p_J} h_{\frac{\phi}{\phi + \bar{Y}_J}}(p_J)$ and $\forall 0 < a < 1$, $h_a(x) \geq \frac{1-x}{1-a} \log^2 \left(\frac{1-x}{1-a} \right)$. Then on $\Omega_{m_f}(\epsilon)$

$$K(\bar{s}_m, \hat{s}_m) \geq \phi \sum_{J \in m} \frac{|J|}{p_J} \frac{1-p_J}{\frac{\bar{Y}_J}{\phi + \bar{Y}_J}} \log^2 \left(\frac{\frac{\bar{Y}_J}{\phi + \bar{Y}_J}}{1-p_J} \right).$$

Introducing

$$V_m^2 = \sum_{J \in m} \phi |J| \frac{1-p_J}{p_J} \log^2 \left(\frac{\frac{\bar{Y}_J}{\phi + \bar{Y}_J}}{1-p_J} \right), \quad (21)$$

we get

$$K(\bar{s}_m, \hat{s}_m) \geq V_m^2, \quad (22)$$

and since $\bar{Y}_J - \phi \frac{1-p_J}{p_J} = \frac{\phi + \bar{Y}_J}{p_J} \left(\frac{\bar{Y}_J}{\phi + \bar{Y}_J} - (1-p_J) \right)$, we have

$$V_m^2 = \sum_{J \in m} |J| \left(\frac{\phi}{\phi + \bar{Y}_J} \right)^2 \frac{\left(\bar{Y}_J - \phi \frac{1-p_J}{p_J} \right)^2}{\phi \frac{1-p_J}{p_J}} \left[\frac{\log \left(\frac{\frac{\bar{Y}_J}{\phi + \bar{Y}_J}}{1-p_J} \right)}{\frac{\frac{\bar{Y}_J}{\phi + \bar{Y}_J}}{1-p_J} - 1} \right]^2.$$

And finally,

$$K(\bar{s}_m, \hat{s}_m) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \geq \rho_{min}^2 \frac{(1-\epsilon)^2}{(1+\epsilon)^4} \chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)}.$$

Moreover, on one hand we have $|m| \leq \mathbf{E}[\chi_m^2] \leq \frac{1}{\rho_{min}} |m|$. On the other hand, since $\chi_m^2 \leq \frac{1}{\Gamma(\log(n))^2 \phi(1-\rho_{max})} (\sum_t Y_t - \sum_t E_t)^2$, using Cauchy-Schwarz Inequality, we get

$$\begin{aligned} \mathbf{E} \left[\chi_m^2 \mathbf{1}_{\Omega_{m_f}(\epsilon)^C} \right] &\leq \frac{\left[\mathbf{E} (Y_t - E_t)^4 + 6\phi^2 \sum_{(t,l), l \neq t} \frac{1-p_t}{p_t^2} \frac{1-p_l}{p_l^2} \right]^{1/2}}{\Gamma(\log(n))^2 \phi(1-\rho_{max})} P(\Omega_{m_f}(\epsilon)^C)^{1/2}, \\ &\leq C(\Gamma, \rho_{min}, \rho_{max}) n^\alpha P(\Omega_{m_f}(\epsilon)^C)^{1/2}, \\ &\leq \frac{C(\phi, \Gamma, \rho_{min}, \rho_{max}, \epsilon, a)}{n^{a/2-\alpha}}, \end{aligned}$$

where $\alpha = 1 - 2 \frac{\log(\log(n))}{\log(n)}$, $n \geq 2$. Finally, we have

$$K(s, \bar{s}_m) + \rho_{min}^2 \frac{(1-\epsilon)^2}{(1+\epsilon)^4} |m| - \frac{C(\phi, \Gamma, \rho_{min}, \rho_{max}, \epsilon, a)}{n^{a/2-\alpha}} \leq \mathbf{E}[K(s, \hat{s}_m)].$$

5.2 Proof of proposition 4.1

Poisson case

The term to be controlled is $\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'}) = \sum_{J \in m'} |J| (\bar{Y}_J - \bar{\lambda}_J) \log \frac{\bar{Y}_J}{\bar{\lambda}_J}$. Using Cauchy-Schwarz inequality, we have

$$\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'}) \leq \sqrt{\chi_{m'}^2} \sqrt{V_{m'}^2},$$

with $\chi_{m'}^2$ and $V_{m'}^2$ defined as in equations (9) and (19). Then, using equation (18)

$$(\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq \sqrt{\chi_{m'}^2} \sqrt{\frac{2}{1-\epsilon} K(\bar{s}_{m'}, \hat{s}_{m'})},$$

and using $2ab \leq \kappa a^2 + \kappa^{-1} b^2$ for all $\kappa > 0$, we get

$$(\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq \frac{\kappa}{2} \chi_{m'}^2 + \frac{\kappa^{-1}}{1-\epsilon} K(\bar{s}_{m'}, \hat{s}_{m'}). \quad (23)$$

And with proposition 2.2, we get, for $\kappa = \frac{1+\epsilon}{1-\epsilon} = 2C(\epsilon)$,

$$\begin{aligned} & (\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon) \cap \Omega_1(\xi)} \\ & \leq \frac{1+\epsilon}{2(1-\epsilon)} \left[|m'| + 8(1+\epsilon) \sqrt{(L_{m'} |m'| + \xi) |m'|} + 4(1+\epsilon)(L_{m'} |m'| + \xi) \right] + \frac{1}{1+\epsilon} K(\bar{s}_{m'}, \hat{s}_{m'}). \end{aligned}$$

Negative binomial case

In this case we can write $\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'}) = \sum_{J \in m'} |J| (\bar{Y}_J - \bar{E}_J) \log \frac{\bar{Y}_J}{\frac{\phi + \bar{Y}_J}{1-p_J}}$. Again, using Cauchy-Schwarz inequality, and with χ_m^2 and V_m^2 defined by equations (9) and (21), we get

$$\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'}) \leq \sqrt{\chi_{m'}^2} \sqrt{V_{m'}^2},$$

so that with equation (22) and $2ab \leq \kappa a^2 + \kappa^{-1} b^2$ for all $\kappa > 0$

$$(\bar{\gamma}(\bar{s}_{m'}) - \bar{\gamma}(\hat{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon)} \leq \frac{\kappa}{2} \chi_{m'}^2 + \frac{\kappa^{-1}}{2} K(\bar{s}_{m'}, \hat{s}_{m'}). \quad (24)$$

Finally, with proposition 2.2 and $\kappa = \frac{1+\epsilon}{2} = 2C(\epsilon)$,

$$\begin{aligned} & (\bar{\gamma}(\hat{s}_{m'}) - \bar{\gamma}(\bar{s}_{m'})) \mathbf{1}_{\Omega_{m_f}(\epsilon) \cap \Omega_1(\xi)} \\ & \leq \frac{1+\epsilon}{4} \left[|m'| + 8(1+\epsilon) \sqrt{(L_{m'} |m'| + \xi) |m'|} + 4(1+\epsilon)(L_{m'} |m'| + \xi) \right] + \frac{1}{1+\epsilon} K(\bar{s}_{m'}, \hat{s}_{m'}). \end{aligned}$$

5.3 Proof of proposition 4.2

Poisson case

Noting that $\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)}] = -\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}]$, we have

$$\begin{aligned}
|\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)}]| &\leq |\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}]| \leq \mathbf{E}[|(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}] \\
&\leq \mathbf{E}\left[\left|\left(\sum_J \sum_t (Y_t - E_t) \log(\rho_{max}/\rho_{min})\right)\right|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}\right] \\
&\leq \log(\rho_{max}/\rho_{min}) \times \mathbf{E}\left[\left|\sum_t (Y_t - E_t)\right|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}\right] \\
&\leq \log(\rho_{max}/\rho_{min}) \times \left(\left[\mathbf{E}\left(\sum_t (Y_t - E_t)^2\right)\right]^{1/2} \times (P(\Omega_{m_f}(\epsilon)^C)^{1/2}\right) \\
&\leq (n\rho_{max})^{1/2} \times \log(\rho_{max}/\rho_{min}) \times (P(\Omega_{m_f}(\epsilon)^C)^{1/2},
\end{aligned}$$

which concludes the proof.

Negative binomial case

Once again, $\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)}] = -\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}]$, and

$$\begin{aligned}
|\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)}]| &\leq |\mathbf{E}[(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}]| \leq \mathbf{E}[|(\bar{\gamma}(\bar{s}_m) - \bar{\gamma}(s))|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}] \\
&\leq \mathbf{E}\left[\left|\left(\sum_J \sum_t \left(Y_t - \phi \frac{1-p_t}{p_t}\right) \log(1/(1-\rho_{min}))\right)\right|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}\right] \\
&\leq \log(1/(1-\rho_{min})) \times \mathbf{E}\left[\left|\sum_t (Y_t - E_t)\right|\mathbf{1}_{\Omega_{m_f}(\epsilon)^C}\right] \\
&\leq \left(n\phi \frac{1}{\rho_{min}^2}\right)^{1/2} \times \log \frac{1}{1-\rho_{min}} \times (P(\Omega_{m_f}(\epsilon)^C)^{1/2}
\end{aligned}$$

which concludes the proof.

5.4 Proof of proposition 4.3

Using the Markov inequality $\mathbf{P} [\bar{\gamma}(s) - \bar{\gamma}(u) \geq b] \leq \inf_a [e^{-ab} \mathbf{E} (e^{a(\bar{\gamma}(s) - \bar{\gamma}(u))})]$ with $a = \frac{1}{2}$, we get

$$\begin{aligned}
\mathbf{P} [\bar{\gamma}(s) - \bar{\gamma}(u) \geq b] &\leq \exp \left[-\frac{b}{2} + \log \mathbf{E} \left[\exp \left(\frac{1}{2} (\gamma(s) - \gamma(u)) + \frac{1}{2} \mathbf{E} [\gamma(u) - \gamma(s)] \right) \right] \right] \\
&\leq \exp \left[-\frac{b}{2} + \frac{1}{2} K(s, u) + \log \mathbf{E} \left[\exp \left(-\frac{1}{2} \sum_t \log \mathbf{P}_s(X_t = Y_t) + \log \mathbf{P}_u(X_t = Y_t) \right) \right] \right] \\
&\leq \exp \left[-\frac{b}{2} + \frac{1}{2} K(s, u) + \sum_t \log \mathbf{E} \sqrt{\frac{\mathbf{P}_u(X_t = Y_t)}{\mathbf{P}_s(X_t = Y_t)}} \right] \\
&\leq \exp \left[-\frac{b}{2} + \frac{1}{2} K(s, u) + \sum_t \mathbf{E} \sqrt{\frac{\mathbf{P}_u(X_t = Y_t)}{\mathbf{P}_s(X_t = Y_t)}} - n \right] \\
&\leq \exp \left[-\frac{b}{2} + \frac{1}{2} K(s, u) - h^2(s, u) \right]
\end{aligned}$$

where $\mathbf{P}_s = \mathbf{P}$ denote the probability under the distribution s . Thus

$$\mathbf{P} [\bar{\gamma}(s) - \bar{\gamma}(u) \geq K(s, u) - 2h^2(s, u) + 2x] \leq e^{-x}.$$

The authors wish to thank Stéphane Robin for more than helpful discussions on the statistical aspect and Gavin Sherlock for his insight on the biological applications.

References

- [1] Risso D, Schwartz K, Sherlock G, Dudoit S: **GC-Content Normalization for RNA-Seq Data**. *BMC Bioinformatics* 2011, **12**:480.
- [2] Braun JV, Muller HG: **Statistical Methods for DNA Sequence Segmentation**. *Statistical Science* 1998, **13**(2):142–162, [<http://www.jstor.org/stable/2676755>].
- [3] Biernacki C, Celeux G, Govaert G: **Assessing a Mixture Model for Clustering with the Integrated Completed Likelihood**. *IEEE Trans. Pattern Anal. Machine Intel.* 2000, **22**(7):719–725.
- [4] Luong TM, Rozenholc Y, Nuel G: **Fast estimation of posterior probabilities in change-point models through a constrained hidden Markov model**. *Arxiv preprint arXiv:1203.4394* under review.
- [5] Akaike H: **Information Theory and Extension of the Maximum Likelihood Principle**. *Second international symposium on information theory* 1973, :267–281.
- [6] Yao YC: **Estimating the number of change-points via Schwarz' criterion**. *Statistics & Probability Letters* 1988, **6**(3):181–189.

- [7] Birgé L, Massart P: **Minimal penalties for Gaussian model selection.** *Probab. Theory Related Fields* 2007, **138**(1-2):33–73.
- [8] Zhang NR, Siegmund DO: **A modified Bayes information criterion with applications to the analysis of comparative genomic hybridization data.** *Biometrics* 2007, **63**:22–32. [PMID: 17447926].
- [9] Braun JV, Braun RK, Muller HG: **Multiple changepoint fitting via quasilielihood, with application to DNA sequence segmentation.** *Biometrika* 2000, **87**(2):301–314, [<http://www.jstor.org/stable/2673465>].
- [10] Birgé L, Massart P: **From model selection to adaptive estimation.** In *Festschrift for Lucien Le Cam*, New York: Springer 1997:55–87.
- [11] Barron A, Birgé L, Massart P: **Risk bounds for model selection via penalization.** *Probab. Theory Related Fields* 1999, **113**(3):301–413.
- [12] Akakpo N: **Estimating a discrete distribution via histogram selection.** *ESAIM Probab. Statist.* 2009, **To appear**.
- [13] Reynaud-Bouret P: **Adaptive estimation of the intensity of inhomogeneous Poisson processes via concentration inequalities.** *Probab. Theory Related Fields* 2003, **126**:103–153.
- [14] Birgé L: **Model selection for Poisson processes.** In *Asymptotics: particles, processes and inverse problems, Volume 55 of IMS Lecture Notes Monogr. Ser.*, Beachwood, OH: Inst. Math. Statist. 2007:32–64.
- [15] Baraud Y, Birgé L: **Estimating the intensity of a random measure by histogram type estimators.** *Probab. Theory Related Fields* 2009, **143**(1-2):239–284.
- [16] Lebarbier E: **Detecting multiple change-points in the mean of Gaussian process by model selection.** *Signal Processing* 2005, **85**(4):717–736.
- [17] Castellan G: **Modified Akaike’s criterion for histogram density estimation.** *C. R. Acad. Sci., Paris, Sér. I, Math.* 330 2000, **8**:729–732.
- [18] Massart P: *Concentration inequalities and model selection, Volume 1896 of Lecture Notes in Mathematics.* Berlin: Springer 2007. [Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard].
- [19] Birgé L, Massart P: **Gaussian model selection.** *J. Eur. Math. Soc. (JEMS)* 2001, **3**(3):203–268.
- [20] Rigai G: **Pruned dynamic programming for optimal multiple change-point detection.** *Arxiv:1004.0887* 2010, [<http://arxiv.org/abs/1004.0887>].
- [21] Cleynen A, Koskas M, Rigai G: **A Generic Implementation of the Pruned Dynamic Programing Algorithm.** *Arxiv preprint arXiv:1204.5564* 2012.

- [22] Arlot S, Massart P: **Data-driven calibration of penalties for least-squares regression.** *J. Mach. Learn. Res.* 2009, **10**:245–279 (electronic), [<http://www.jmlr.org/papers/volume10/arlot09a/arlot09a.pdf> [pdf]].
- [23] Johnson N, Kemp A, Kotz S: **Univariate Discrete Distributions.** *John Wiley & Sons, Inc.* 2005.
- [24] Killick R, Eckley I: **changepoint: An R package for changepoint analysis** 2011.
- [25] Breiman L, Friedman J, Olshen R, Stone C: *Classification and Regression Trees.* Monterey, CA: Wadsworth and Brooks 1984.
- [26] Rigaiil G, Lebarbier E, Robin S: **Exact posterior distributions and model selection criteria for multiple change-point detection problems.** *Statistics and Computing* 2012, **22**:917–929.